

(51) International Patent Classification ⁶ : C12Q 1/68, C07H 21/04, C12N 15/00	A1	(11) International Publication Number: WO 99/22025 (43) International Publication Date: 6 May 1999 (06.05.99)
(21) International Application Number: PCT/US98/22519 (22) International Filing Date: 23 October 1998 (23.10.98) (30) Priority Data: 60/063,103 24 October 1997 (24.10.97) US Not furnished 22 October 1998 (22.10.98) US (71) Applicant: TRUSTEES OF BOSTON UNIVERSITY [US/US]; 147 Bay State Road, Boston, MA 02215 (US). (72) Inventors: CANTOR, Charles, R.; 11 Bay State Road #6, Boston, MA 02115 (US). SIDDIQI, Fouad, A.; Apartment #3, 113 Bay State Road, Boston, MA 02115 (US). (74) Agents: CARROLL, Peter, G. et al.; Medlen & Carroll, LLP, Suite 2200, 220 Montgomery Street, San Francisco, CA 94104 (US).		(81) Designated States: AU, CA, JP, European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE). Published <i>With international search report.</i>
(54) Title: INFERENCE SEQUENCING BY HYBRIDIZATION (57) Abstract <p>The inventors present a new strategy for <i>de novo</i> sequencing using oligonucleotide probe arrays which we term Inference Sequencing by Hybridization (ISBH). The ISBH method characterises a target DNA with several small degenerate probe arrays and reconstructs the base sequence of the target via computer algorithm at 100 % accuracy. ISBH returns a target sequence as multiple non-overlapping fragments several hundred to several thousand bases in length. Target DNAs up to 100 kilobases long can potentially be sequenced in a single experiment. A hypothetical ISBH probe array and sequence reconstruction algorithm were designed and tested by computer simulation on 76 DNA sequences obtained from GenBank comprising a total of 2.45 million base pairs. ISBH performs optimally on target DNAs in the range of 15-45 kilobases that have low levels of repeated sequence. We believe that the ISBH approach has great potential to increase the speed of DNA sequencing.</p>		

PCT

WORLD INTELLECTUAL PROPERTY ORGANIZATION
International Bureau



FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakhstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

INFERENCE SEQUENCING BY HYBRIDIZATION

This application for patent under 35 U.S.C. § 111(a) claims priority to Provisional Application Serial No. 60/063,103, filed October 24, 1997 under 35 U.S.C. § 111(b). This invention was made with Government Support under Contract Number DGE-9452651 awarded by the National Science Foundation. The Government has
5 certain rights in the invention.

BACKGROUND OF THE INVENTION

Procedures involving use of Sequencing by hybridization (SBH) are known to those skilled in the art, and have recently been demonstrated to be useful as a powerful
10 alternative to electrophoretic methods for diagnostic DNA analysis [M. Chee *et al.*, Accessing Genetic Information With High Density DNA Arrays, *Science* 274, 610 (1996); J. Hacia *et al.*, Detection Of Heterozygous Mutations In BRCA1 Using High-Density Oligonucleotide Arrays And Two-Color Fluorescence Analysis, *Nature Genetics* 14, 441 (1996)]. Diagnostic SBH employs hybridization of a target DNA
15 sequence to a tiled array of several thousand short oligonucleotide probes of known sequence [W. Bains *et al.*, A Novel Method For Nucleic Acid Sequence Determination, *J Theor Biol*, 135, 303 (1988); E. Southern *et al.*, Hybridization With Oligonucleotide Arrays, *Genomics*, 13, 1008 (1992)]. The pattern of hybridization, detected by fluorescence microscopy, indicates which oligonucleotides in the probe
20 array are present in the target DNA. When this information is compared against a reference target sequence, the entire sequence of the target DNA can be reconstructed at high accuracy.

SBH, while offering great advantages in terms of throughput for diagnostic sequence analysis, suffers from the drawback that a different probe array must be
25 tailored for each target DNA analyzed [M. Chee *et al.*, Accessing Genetic Information With High Density DNA Arrays, *Science*, 274, 610 (1996)]. For *de novo* sequencing, current SBH methods are not competitive with electrophoretic sequencing techniques

that yield 600-1000 base pair read lengths per experiment [P. Pevzner *et al.*, Improved Chips For Sequencing By Hybridization, *J Biomolecular Struct Dyn*, 9, 399 (1991)].

Even under perfect experimental conditions, existing SBH designs cannot reconstruct a unique target sequence from hybridization data alone [P. Pevzner *et al.*, Towards DNA

5 Sequencing Chips, In *19th Int Conf Mathematical Foundations of Computer Science*, Lecture Notes in Computer Science, Springer-Verlag, Berlin, Vol 841, pp. 143-158

(1994); P. Pevzner, Rearrangements Of DNA Sequences And SBH, *Computers Chem*, 18, 221 (1994)]. Without a reference sequence for comparison, *de novo* SBH is

10 fundamentally limited because it acquires base sequence information at the cost of positional information. One knows exactly *which* subsequences (probe sequences) are present in the target DNA but not *where* they are located.

Subsequences must be arranged by examining how they overlap with one another. For example, octa-nucleotides are assembled into longer sequences by finding corresponding seven-base overlaps. The accuracy of reassembly is limited because any

15 particular subsequence can occur more than once in the target DNA, leading to an ambiguity in the final reconstructed sequence [P. Pevzner *et al.*, Towards DNA Sequencing Chips, In *19th Int Conf Mathematical Foundations of Computer Science*, Lecture Notes in Computer Science, Springer-Verlag, Berlin, Vol 841, pp. 143-158

(1994); P. Pevzner, Rearrangements Of DNA Sequences And SBH, *Computers Chem*, 20 18, 221 (1994)]. *De novo* SBH designs typically use the complete set of all possible oligonucleotide probes of a given length [W. Bains *et al.*, A Novel Method For Nucleic Acid Sequence Determination, *J Theor Biol*, 135, 303 (1988); N. Broude *et al.*, Enhanced DNA Sequencing By Hybridization, *Proc Natl Acad Sci USA*, 91, 3072

(1994); R. Drmanac *et al.*, DNA Sequence Determination by Hybridization: A

25 Strategy For Efficient Large-Scale Sequencing, *Science*, 260, 1649 (1993); R. Drmanac *et al.*, Sequencing Of Megabase Plus DNA By Hybridization: Theory Of The Method, *Genomics*, 4, 114 (1989)]; the use of longer probes increases reconstruction accuracy but requires the use of very large arrays ($> 10^9$ probes), since the number of required probes increases exponentially with probe length. Even assuming perfect hybridization,

an SBH array containing all $\sim 10^6$ possible 10-mers would reliably be able to sequence only about 600 bp of target DNA in a single experiment. As longer target DNA sequences are attempted, the reconstruction accuracy drops precipitously. For a detailed discussion of SBH reassembly algorithms and their limitations see references [P. Pevzner *et al.*, Towards DNA Sequencing Chips, In *19th Int Conf Mathematical Foundations of Computer Science*, Lecture Notes in Computer Science, Springer-Verlag, Berlin, Vol 841, pp. 143-158 (1994); P. Pevzner *et al.*, Improved Chips For Sequencing By Hybridization, *J Biomolecular Struct Dyn*, 9, 399 (1991); P. Pevzner, Rearrangements Of DNA Sequences And SBH, *Computers Chem*, 18, 221 (1994)].

All *de novo* SBH strategies proposed thus far are direct methods, that is, they directly probe the target DNA with oligonucleotides whose sequences are then assembled into longer fragments. Since it is currently not feasible to manufacture a probe array containing more than approximately 10^6 tiled oligonucleotide probes, a direct *de novo* SBH approach cannot outperform electrophoretic sequencing in terms of read length and reaccuracy.

IN THE FIGURES

The present invention will become better understood with reference to the following description, appended claims, and accompanying figures where:

Figure 1. Basic concepts underlying the design of an Inference Sequencing by Hybridization (ISBH) probe array. A target, in this case a four-digit phone number, is characterized by a degenerate probe array with a four-fold redundancy using only 64 probes. The probe array does not detect any digit directly, but the information gathered is sufficient to unambiguously infer the identity of the number. A probe array based on conventional SBH designs capable of acquiring the same information would require 10,000 probes, one for each phone number.

Figure 2. General scheme for ISBH. A long single-stranded target DNA is sheared into short oligonucleotides and hybridized to an ISBH probe array. The pattern of hybridization is used to create a set of degenerate 16-mers that characterize

the target DNA. Information from this degenerate set is used by an inference algorithm to produce a set of explicit 16-mers. The set of explicit 16-mers produced by the inference algorithm contains all 16-mers actually present in the target DNA sequence as well as "false positive" 16-mers that are not in the target. A data reduction algorithm is then used to eliminate the false positives from the set of explicit 16-mers. The explicit 16-mers that remain after data reduction are then reassembled at high accuracy into contiguous sequence.

Figure 3. Design of the ISBH probe array of degenerate 16-mers used in this study. The array consists of 25 different probe groups. Each group pattern represents $2^{16} = 65,536$ degenerate 16-mers, for a total of 1,638,400 probes. Each probe in the array represents 65,536 explicit 16-mers. Under ideal conditions, a single target 16-mer will hybridize to exactly one probe in each probe group.

Figure 4. Example of how false positives are generated by inference. Tetranucleotides from a target DNA are characterized by an ISBH probe array of degenerate 4-mers with two-fold redundancy ($R = A$ or G ; $Y = C$ or T ; $W = A$ or T ; $S = C$ or G). The inference algorithm generates all valid combinations of data from the probe array to produce a set of nine inferred tetranucleotides, six of which are false positives. In general, the number of false positives generated by inference decreases with the number of probes used in the ISBH array. If the ISBH array used in this example had 16 additional probes, then no false positives would have been generated.

Figure 5. Data reduction after inference for the ISBH array of degenerate 16-mers used in this study. Seventy-six target DNA sequences downloaded from GenBank comprising a total of 2.45 million bases were tested by computer simulation. The number of 16-mers generated by inference increases as a power law function of the number of different 16-mers in the target DNA (filled circles). Data reduction reliably eliminated all but a handful of the false positives for all target lengths investigated (open triangles), even when false positives comprised more than 99% of the inferred 16-mer set.

Figure 6. Absolute performance of the ISBH sequence reassembly algorithm. The reassembly algorithm returns a target DNA as several non-overlapping fragments

in the range of several hundred to several thousand bases in length. The largest fragment reconstructed in this study was 28 kilobases, and fragments longer than 10kb were commonly observed. Reconstructed fragments always show 100% identity to some region of the target DNA.

5 Figure 7. Total target coverage in a simulated ISBH experiment. ISBH typically covers more than 95% of a target DNA in a single hybridization experiment in fragments that are longer than 500 bases.

 Figure 8. Summary of ISBH simulation results. All lengths are in bases. Locus and definition for each sequence are shown exactly as they appear in GenBank.
10 The number different 16-mers in a target is defined as the number of 16-mers which have different base sequences. The fraction of target 16-mers that are repeated is given by: $1 - [(\text{number of different 16-mers in target})/(\text{target length} - 15)]$, which is a quantitative measure of repetitiveness of the sequence. The fraction of the inferred set that are false positive is given by: $1 - [(\text{number of different 16-mers in target})/(\text{number of 16-mers in inferred set})]$.
15

SUMMARY OF THE INVENTION

 The present invention is directed to a method that satisfies the above mentioned problems by introducing a new SBH implementation for *de novo* sequencing, which we term Inference Sequencing by Hybridization (ISBH). The basic concepts
20 underlying ISBH are illustrated in Figure 1 with an example of how to determine the last four digits of a phone number without detecting any digit directly. A conventional SBH strategy would require $10^4 = 10,000$ probes (one for each phone number), whereas the ISBH approach gathers the same information using only $4 \times 2^4 = 64$ degenerate probes. The digit groupings in Figure 2 are analogous to familiar base
25 groupings such as purine/pyrimidine (R/Y), amino/keto (M/K), and weak/strong (W/S).

 ISBH is an indirect strategy that uses several small arrays (65,536 probes each) to closely approximate the information that would be gathered from a single SBH array containing ~4.3 billion probes (Figure 2). Our strategy relies on degenerate probe arrays that are similar to binary SBH arrays proposed by Pevzner *et. al.* [P. Pevzner *et*

al., Towards DNA Sequencing Chips, In *19th Int Conf Mathematical Foundations of Computer Science*, Lecture Notes in Computer Science, Springer-Verlag, Berlin, Vol 841, pp. 143-158 (1994)]. Unlike conventional SBH, whose accuracy drops with increasing target length, we find that ISBH always reconstructs targets at 100% accuracy. ISBH returns a target sequence as non-overlapping fragments that range from several hundred to several thousand bases in length. The number of fragments increases with the length and repetitiveness of the target DNA, but in principle, any length of target can be sequenced in a single experiment. We have designed and empirically tested by computer simulation an ISBH array and reconstruction algorithm for use in *de novo* sequencing of targets up to 100 kilobases in length.

METHODS

The inventors have simulated the following laboratory experiment: 1) a single stranded target DNA of unknown sequence is sheared into overlapping oligomers 16 bases long. 2) These oligomers are hybridized to an ISBH degenerate probe array and this pattern of hybridization is detected. 3) the hybridization data is reconstructed by computer algorithm into contiguous sequence.

An ISBH probe array containing 1.64 million different degenerate 16-mers was designed. This array is 25-fold redundant - it consists of 25 groups of $2^{16}=65,536$ (64K) different probe oligomers. Figure 3 shows the identity of each degenerate probe group in the array. Each group of 64K degenerate probes is capable of hybridizing to all possible explicit 16-mers. Thus, a single explicit 16-mer will, under ideal conditions, hybridize to exactly 25 different degenerate probes in the ISBH array.

Simulations of ISBH experiments and subsequent data analysis were performed on a Silicon Graphics Origin2000 supercomputer (Boston University Center for Computational Science) with code written in C as follows: 1) Each target DNA sequence was retrieved from GenBank ([http:// www2.ncbi.nlm.nih.gov/genbank](http://www2.ncbi.nlm.nih.gov/genbank)) and broken up into all possible component 16-mers. This set of 16-mers was then shuffled to destroy any positional information. 2) Each 16-mer from the target DNA was compared against all 1.64 million degenerate probes in the ISBH array, simulating

ideal hybridization. If a 16-mer from the target was found to hybridize with one of the degenerate probes, then the hybridization was considered to be a signal arising from that particular probe in the ISBH array. 3) The 16-mers from the target DNA were then discarded, to account for the fact that in a real experiment, the target DNA is of unknown sequence. 4) The signals from the ISBH array were collected to produce a set of degenerate 16-mers, which is the non-repeated set of degenerate 16-mers present in the target DNA. We emphasize that the only data retained is the signal pattern from the ISBH array, as would be detected in an actual ISBH experiment. The degenerate set contains no information about the explicit base sequence of any 16-mer present or its position in the target DNA sequence.

Sequence Reconstruction

1) *Inference.* A set of explicit 16-mers is inferred from the set of degenerate 16-mers detected by the ISBH array. The inference is accomplished by testing every possible explicit 16-mer against the degenerate set. Any particular explicit 16-mer is included in the inferred set only if exactly 25 corresponding degenerate 16-mers are present in the data from the ISBH array. Under conditions of ideal hybridization, the inferred 16-mer set is always a superset of the set of 16-mers present in the target DNA sequence. The inferred set usually contains false positive 16-mers, ones which are not actually present in the original target DNA sequence. The number of these false positives increases with the length and repetitiveness of the target DNA (Figure 4).

2) *Primary Data Reduction.* Since the inferred set of explicit 16-mers contains an unknown number of false positives, a data reduction step is required to eliminate them. Every 16-mer in the inferred set is examined to determine if it overlaps by six bases at least one other 16-mer in the inferred set on both its 3' and 5' ends. If both overlaps are not found the 16-mer is discarded. This procedure is repeated iteratively on the resultant set of 16-mers, each time examining an overlap one base longer than was used for the previous iteration, until an overlap of fifteen bases is reached. The

set that remains is then iterated using fifteen-base overlaps until four or fewer 16-mers are discarded at each iteration.

5 3) *Secondary Data Reduction.* All possible reconstruction ambiguities are eliminated by comparing the 3'-fifteen bases of each 16-mer with the 3'-fifteen bases of every other 16-mer in the set from step 2 above. If two or more 16-mers are found to have identical fifteen base 3'-ends, then they are all discarded from the data set. A similar procedure is used to compare the 5'-fifteen bases of the 16-mers and eliminate any duplication.

10 4) *Sequence Reassembly.* The 16-mers remaining in the inferred set are then assembled into longer sequences. This is done by comparing the 3'-fifteen bases of each 16-mer to the 5'-fifteen bases of every other 16-mer in the data set. If a match is found, then the two 16-mers are combined into a single 17-mer. The two terminal 15-mers on either side of this newly formed 17-mer are compared for overlap with the remaining 16-mers in the data set and the process is repeated until no more overlaps
15 are found. To insure reconstruction accuracy, only fragments at least 100 bases in length were considered to be part of the target sequence.

RESULTS

Inference Algorithm

20 ISBH simulations were performed on 76 target sequences obtained from GenBank comprising a total of 2.45 million bases, ranging from 5 to 100 kilobases in length. The size of the inferred 16-mer pool increases as a power law function of the number of different 16-mers in the target DNA. Data reduction reliably eliminates all but a handful of the false positives for all target lengths investigated, even when false
25 positives accounted for more than 99% of the inferred 16-mer pool (Figure 5). The set of inferred 16-mers remaining after data reduction closely approximates the information that would be gathered from an SBH array containing all explicit 16-mers (~4.3 billion probes).

Sequence Reconstruction

The ISBH reconstruction algorithm typically returns a target DNA as several non-overlapping fragments that are in the range of several hundred to several thousand bases in length. Reconstructed fragments always show 100% identity (by BLASTN) to some region of unknown position in the target DNA sequence. The largest single fragment reconstructed was 28 kb, and fragments longer than 10kb were commonly observed (Figure 6). In most cases, more than 95% of the target DNA was recovered in a simulated ISBH experiment (Figure 7). ISBH performs optimally on target DNAs having no repeated 16-mers, generally returning a handful of long (3-15 kb) fragments. Even on sequences with many repeated 16-mers, ISBH returns dozens of fragments shorter than 5 kb, which is five to ten times the performance of electrophoretic sequencing methods. Target DNAs longer than 50 kb tend to produce large numbers of false positives in the inference step, a few of which remain after data reduction. False positives introduce ambiguities during reassembly, leading to lower average lengths for reconstructed fragments. A comprehensive list of each sequence analyzed as well as a summary of the ISBH simulation results are shown in Figure 8.

DETAILED DESCRIPTION OF THE INVENTION

While this invention is satisfied by embodiments in many different forms, there will herein be described preferred embodiments of the invention, with the understanding that the present disclosure is to be considered exemplary of the principles of the invention and is not intended to limit the invention to the embodiments illustrated and described. The scope of the invention will be measured by the appended claims and their equivalents.

A benefit of the present invention, in contrast to conventional SBH, is that the inventor's ISBH method has the potential to sequence very long targets at high accuracy, using an oligonucleotide array of moderate size. The hypothetical ISBH array studied here could easily sequence 15-45kb of DNA in a single experiment. The ISBH method requires no electrophoresis, no information about the target DNA, and could be used for diagnostic as well as *de novo* applications. In a single experiment,

ISBH could generate more sequence data than two dozen Sanger sequencing reactions after shotgun subcloning of a target DNA. In the best cases, each fragment reconstructed by the ISBH method can outperform electrophoretic methods by 28-fold. In the worst cases, ISBH performance is equivalent to electrophoretic methods. ISBH reconstruction of a single target DNA generally required less than 10 minutes of supercomputer time to complete. The computational complexity of the inference step is of order N^2 , while the data reduction and reassembly steps are of order $N \log(N)$. Sequence reconstruction using a highly streamlined ISBH algorithm running on a typical desktop computer could be completed in a few hours.

For DNA of random sequence, a given 16-mer should appear once every $4^{16} \approx 4.3 \times 10^9$ bases, a 15-mer once in $4^{15} \approx 10^9$ bases, and a 14-mer once in $4^{14} \approx 2.7 \times 10^8$ bases. We note however, that DNAs from a wide variety of organisms in the range of 10-100 kb typically have hundreds or thousands of repeated 16-mers. As shorter subsequences are examined, the number of repeated subsequences increases dramatically. For example, the 48.5 kb genome of bacteriophage lambda, which has no repeated 16-mers, has a single repeated 15-mer, and ten distinct 14-mers appearing more than once. This would suggest that any form of *de novo* SBH using oligonucleotide probes shorter than 16 bases will perform poorly on target DNAs longer than a few kilobases. ISBH appears to perform optimally both in terms of absolute read lengths and relative target coverage on DNAs in the range of 25kb that have small numbers of repeated 16-mers. ISBH is a scalable technique - the number of false positives generated by inference increases as the number of probes used in the ISBH array decreases. An ISBH array smaller than the one examined here (e.g., 12 probe groups using $12 \times 2^{16} = 7.86 \times 10^5$ probes) would still sequence with 100% accuracy, but would return a target as shorter fragments.

While ISBH under ideal conditions would appear to provide an enormous gain for *de novo* sequencing over conventional SBH and electrophoretic methods, several daunting technical obstacles remain. Each degenerate probe is actually a mixture of many individual probes that are bound to the same area in an SBH array. The binding capacity of such a degenerate probe is greatly reduced in comparison to a pure

individual probe - for the hypothetical ISBH array in this study, the complexity of each degenerate probe is 65,536. Such a high probe complexity may mean that accurate physical hybridization cannot be achieved with a high signal to noise ratio. Possible solutions to this problem include the use of base analogs to decrease probe complexity or the addition of an enzymatic step (*e.g.*, ligation) to augment the accuracy of simple physical hybridization.

Noise contamination of the data set, particularly in terms of false negatives, must be studied in greater detail. False positives are easily dealt with in the data reduction step, but false negatives (target 16-mers that never appear at all in the inferred data set) will have the effect of lowering the mean fragment length during reconstruction. Aberrant hybridization also increases the complexity of data processing needed for reliable sequence reconstruction; the upper bound for robust sequence reconstruction from an actual ISBH implementation is likely to be somewhat lower than the ideal situation presented here.

An ISBH sequencing approach would be very effective for rapid analysis of viral and bacterial genomes which are essentially non-repeating. Sequencing of double-stranded target DNAs by ISBH is also possible, as is sequencing of a mixture of targets. For double-stranded DNA, ISBH performance is equivalent to the case of a single-stranded DNA twice as long. If a double stranded target is cleaved by a restriction endonuclease before hybridization, ISBH will return the sequence of each restriction fragment. If this experiment is repeated using a restriction endonuclease with a different recognition site, then the fragments can be aligned relative to one another using standard contig reassembly algorithms. The potential of the ISBH strategy is so strong that we are now investigating strategies to implement it in practice.

Another embodiment of the ISBH Probe Array, Experimental, and Data Analysis Algorithm Design consists of the following:

Probe Array Design. The proposed array consists of 768K oligonucleotide probes divided into three groups: 1) all 256K possible 9-base single-stranded

sequences, 2) all 256K possible 9-base 5'-overhanging partial duplexes, and 3) all 256K possible 9-base 3'-overhanging partial duplexes.

Target Preparation. The target DNA must be single stranded - it may be prepared (for example) by long PCR of a double stranded target using one primer that is biotinylated at its 5' end to facilitate purification from the other strand. The biotinylated strand may be captured on streptavidin-coated beads or column.

Experimental Design. Eight distinct oligonucleotides in two groups are required as follows. Group I: 5'-ANNNNNNN-3', 5'-CNNNNNNN-3', 5'-GNNNNNNN-3', 5'-TNNNNNNN-3'. Group II: 5'-P-NNNNNNNA-Sdd-3', 5'-P-NNNNNNNC-Sdd-3', 5'-P-NNNNNNNG-Sdd-3', and 5'-P-NNNNNNNT-Sdd-3'. P denotes a phosphate group, Sdd denotes a 3'-dideoxy base connected by a phosphorothioate linkage. Sixteen separate reactions are then performed: each oligonucleotide from group I is combined with an oligonucleotide from group II and then hybridized to the single-stranded target under conditions favoring accurate base-pairing. After hybridization has occurred, the oligonucleotides that are still remaining in solution (for example) must be removed by size-exclusion column chromatography. DNA ligase (and all necessary cofactors) are added to the reaction mixture. After the ligation reaction, exonuclease III is added to the reaction mixture to destroy any unligated oligonucleotides that are still hybridized to the target.

Under ideal conditions, the ligation products will be 16-mers of the form: 5'-ANNNNNNNNNNNNNNA-Sdd-3' (and all 15 other permutations of the end bases). The ligation products from each of the sixteen reactions are then hybridized separately to a probe array as described above.

Data Analysis. Each of the sixteen probe array hybridization experiments described above generates the following data: a set of 9-mers from probe group 1, a set of 9-mers from probe group 2, and a set of 9-mers from group 3. The following algorithm is used to expand the data: compare each 9-mer from group 2 to each 9-mers from group 3. If the 9-mer from group 2 has the same 3' two bases as the 5' two bases of the 9-mer in group 3 (i.e., they have a two-base overlap), then combine the two 9-mers to form a single 16-mer (concatenate the 3' seven bases of the group 3

oligo to the 3' end of the group 2 oligo). This newly formed 16-mer is retained only if all eight of its 9-base subsequences are present in group 1. This analysis is performed for all probe array experiments and all retained 16-mers are collected into a single set. This set of 16-mers (which is the inferred set of explicit 16-mers from the target) is then subjected to the same data reduction and sequence reconstruction algorithms that we have previously described for ISBH.

Accordingly, this invention is not limited to the particular embodiments disclosed, but is intended to cover all modifications that are within the spirit and scope of the invention as defined by the appended claims.

CLAIMS

What is claimed is:

1. A method of testing a nucleic acid target, comprising:
 - a) fragmenting a single-stranded target DNA into single-stranded target DNA fragments;
 - b) testing said fragments so as to generate a signal for each fragment;
 - c) calculating a first set of N-mers from said signals, each of said N-mers having a sequence with 3' and 5' ends;
 - d) comparing a portion of the nucleic acid sequence of each of said N-mers of said first set with a portion of the nucleic acid sequence of every other N-mer in said first set for sequence overlap; and
 - e) eliminating each N-mer that is found not to display said overlap, so as to create a second set of N-mers.
2. The method of Claim 1, further comprising the steps:
 - f) comparing a portion of the nucleic acid sequence of each of said N-mers in said second set with a portion of the nucleic acid sequence of every other N-mer in said second set for sequence overlap, wherein the portion compared has a length in bases defined by N-1; and
 - g) identifying an instance where said portion of one N-mer from said second set is found to overlap with said portion of another N-mer from said second set, thereby identifying first and second overlapping N-mers.
3. The method of Claim 1, wherein non-target DNA fragments are added to said single-stranded target DNA fragments prior to step (b).

4. The method of Claim 2, wherein N is 16 and said portion compared in step (f) is fifteen bases in length.

5. The method of Claim 4, wherein said fifteen bases compared in step (f) are the 3'-fifteen bases of each 16-mer with the 5'-fifteen bases of every other 16-mer in said second set.

6. The method of Claim 5, further comprising the step (h) constructing a 17-mer from the combination of the sequences of said first and second overlapping N-mers of step (g).

7. A method of testing a nucleic acid target, comprising:

a) fragmenting a single-stranded target DNA into single-stranded fragments;

b) testing each of said fragments so as to generate a signal for each fragment;

c) calculating a first set of N-mers from said signals, each of said N-mers having a sequence with 3' and 5' ends;

d) comparing a portion of the nucleic acid sequence of each of said N-mers of said first set with a portion of the nucleic acid sequence of every other N-mer in said first set for sequence overlap;

e) eliminating each N-mer that is found not to display said overlap, so as to create a second set of N-mers;

f) comparing a portion of the nucleic acid sequence of each of said N-mers in said second set with a portion of the nucleic acid sequence of every other N-mer in said second set for sequence overlap, wherein the portion compared has a length in bases defined by N-1; and

g) identifying an instance where said portion of one N-mer from said second set is found to overlap with said portion of another N-mer from said second set, thereby identifying first and second overlapping N-mers.

8. The method of Claim 7, wherein non-target DNA fragments are added to said single-stranded target DNA fragments prior to step (b).

9. The method of Claim 8, wherein N is 16 and said portion compared in step (f) is fifteen bases in length.

5 10. The method of Claim 9, wherein said fifteen bases compared in step (f) are the 3'-fifteen bases of each 16-mer with the 5'-fifteen bases of every other 16-mer in said second set.

10 11. The method of Claim 10, further comprising the step (h) constructing a 17-mer from the combination of the sequences of said first and second overlapping N-mers of step (g).

12. A method of testing a nucleic acid target, comprising:

- 15 a) fragmenting a single-stranded target DNA into single-stranded fragments;
- b) adding non-target DNA fragments to said single-stranded target DNA fragments to create a mixture;
- c) testing each of said fragments so as to generate a signal for each fragment;
- d) calculating a first set of N-mers from said signals, each of said N-mers having a sequence with 3' and 5' ends;
- 20 e) comparing a portion of the nucleic acid sequence of each of said N-mers of said first set with a portion of the nucleic acid sequence of every other N-mer in said first set for sequence overlap;
- f) eliminating each N-mer that is found not to display said overlap, so as to create a second set of N-mers;

g) comparing a portion of the nucleic acid sequence of each of said N-mers in said second set with a portion of the nucleic acid sequence of every other N-mer in said second set for sequence overlap, wherein the portion compared has a length in bases defined by N-1; and

5 h) identifying an instance where said portion of one N-mer from said second set is found to overlap with said portion of another N-mer from said second set, thereby identifying first and second overlapping N-mers.

13. The method of Claim 12, wherein N is 16 and said portion compared in step (g) is fifteen bases in length.

10 14. The method of Claim 13, wherein said fifteen bases compared in step (g) are the 3'-fifteen bases of each 16-mer with the 5'-fifteen bases of every other 16-mer in said second set.

15 15. The method of Claim 14, further comprising the step (i) constructing a 17-mer from the combination of the sequences of said first and second overlapping N-mers of step (h).

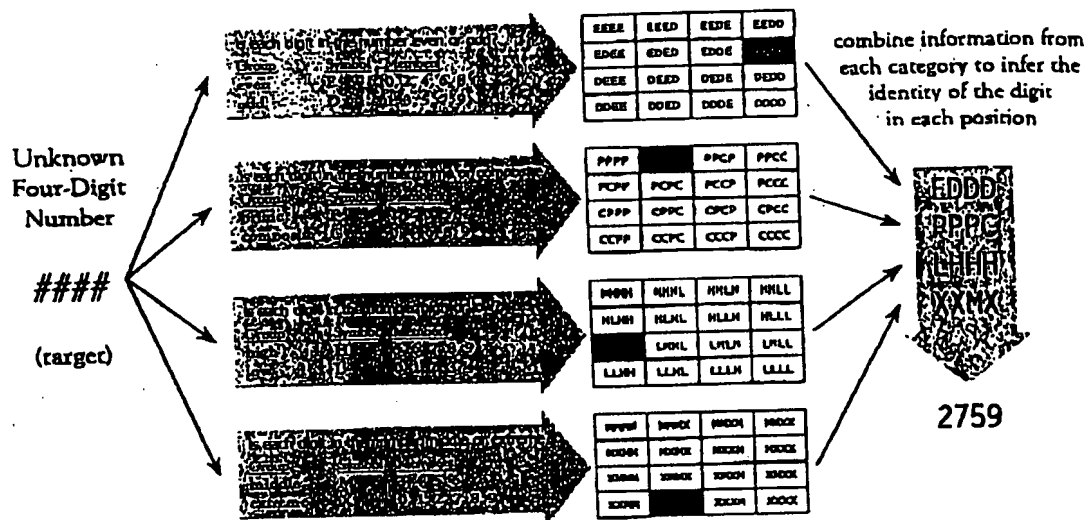


Figure 1

2/11

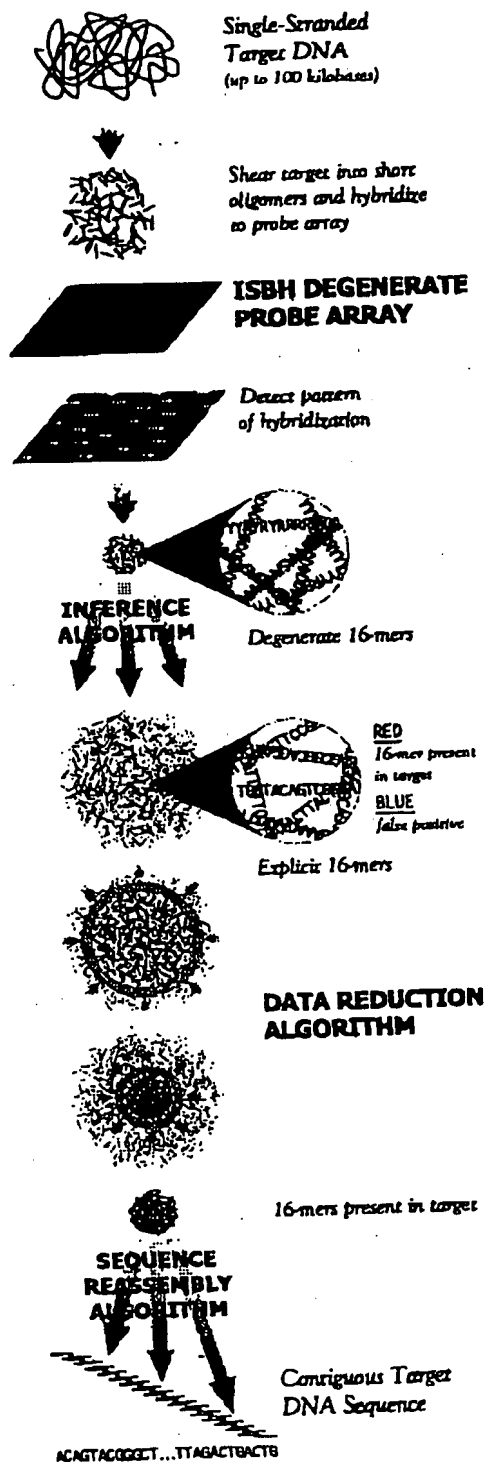


Figure 2

Degenerate Probe Group Pattern	Example Probe
[A/B] ₁₆	ABAABBABABABABBA
[C/D] ₁₆	DCDCDDCCDCDCDCDC
[G/H] ₁₆	GHHHGGGGHGHGHGHG
[T/V] ₁₆	TVTTTTTVVVTTT
[R/Y] ₁₆	RYRRYYRRYYRRYYRR
[M/K] ₁₆	MKKKMMKMMMCKMKM
[W/S] ₁₆	SSWSWSSWSWSWSWS
[R/Y] ₈ [M/K] ₈	RRRRRRRRMMMMMM
[M/K] ₈ [R/Y] ₈	MMMMMMRRRRRRRR
[M/K] ₈ [W/S] ₈	MMMMMMWWWWWW
[W/S] ₈ [M/K] ₈	WWWWMMMMMMMM
[R/Y] ₈ [W/S] ₈	RRRRRRRRWWWWWW
[W/S] ₈ [R/Y] ₈	WWWWRRRRRRRR
[M/K] ₄ [R/Y] ₈ [M/K] ₄	MMMRRRRRRRMM
[W/S] ₄ [R/Y] ₈ [W/S] ₄	WWWRRRRRRWWW
[W/S] ₄ [M/K] ₈ [W/S] ₄	WWWMMMMMMWWW
[R/Y] ₄ [M/K] ₈ [R/Y] ₄	RRRRMMMMMMRRRR
[R/Y] ₄ [W/S] ₈ [R/Y] ₄	RRRRWWWWWWRRRR
[M/K] ₄ [W/S] ₈ [M/K] ₄	MMMMWWWWWWMM
[(R/Y)[M/K]] ₈	RMRMRMRMRMRM
[(W/S)[M/K]] ₈	WMWMWMWMWMWM
[(R/Y)[W/S]] ₈	RWRWRWRWRWRW
[(R/Y) ₄ (M/K) ₄] ₂	RRRRMMMMRRRRMM
[(R/Y) ₄ (W/S) ₄] ₂	RRRRWWRRRRWW
[(W/S) ₄ (M/K) ₄] ₂	WWWWMMMMWWWWMM

R = A or G
Y = C or T

M = A or C
K = G or T

W = A or T
S = C or G

B = C or G or T
D = A or G or T
H = A or C or T
V = A or C or G

Figure 3

4/11

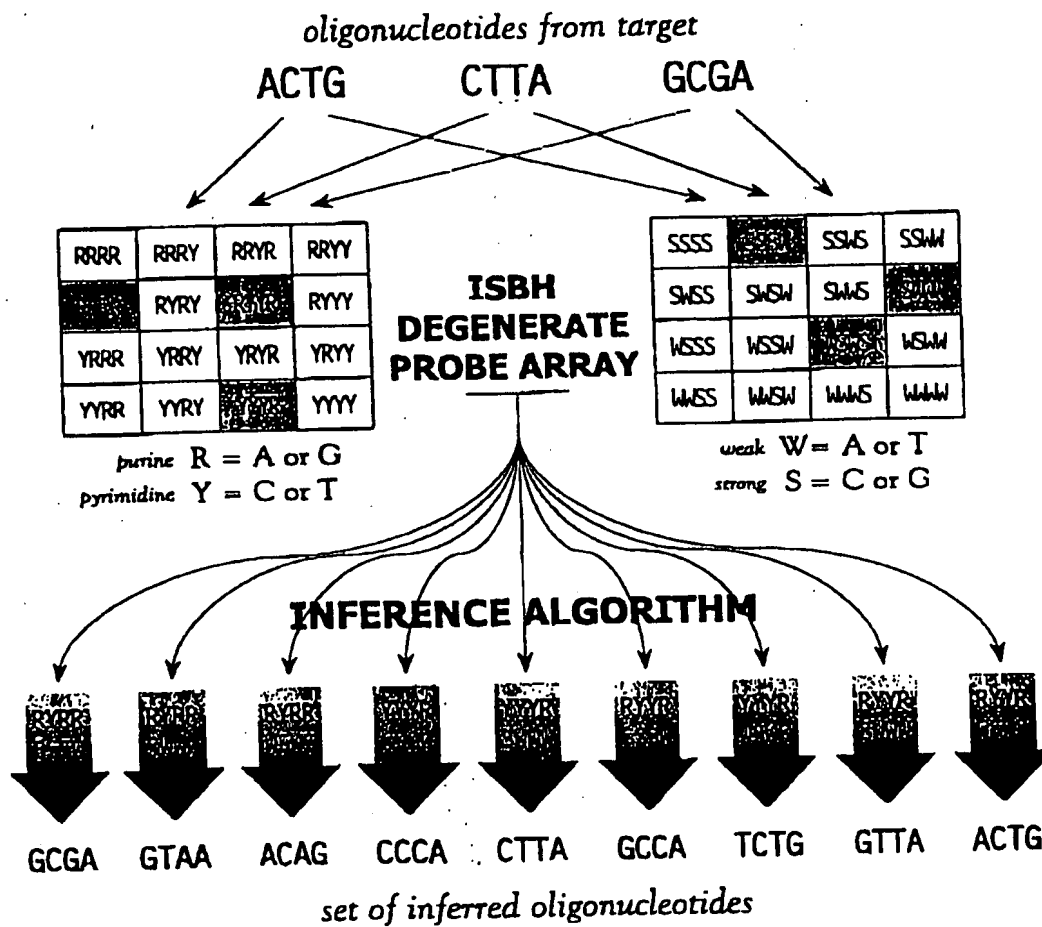


Figure 4

5/11

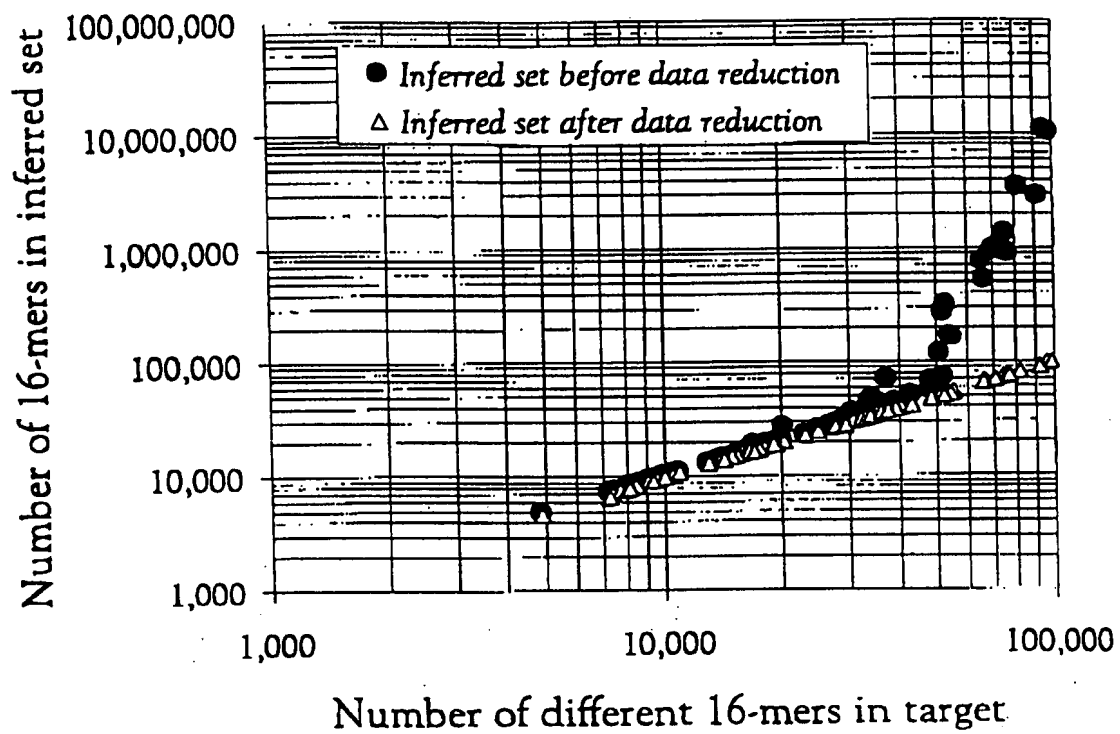


Figure 5

6/11

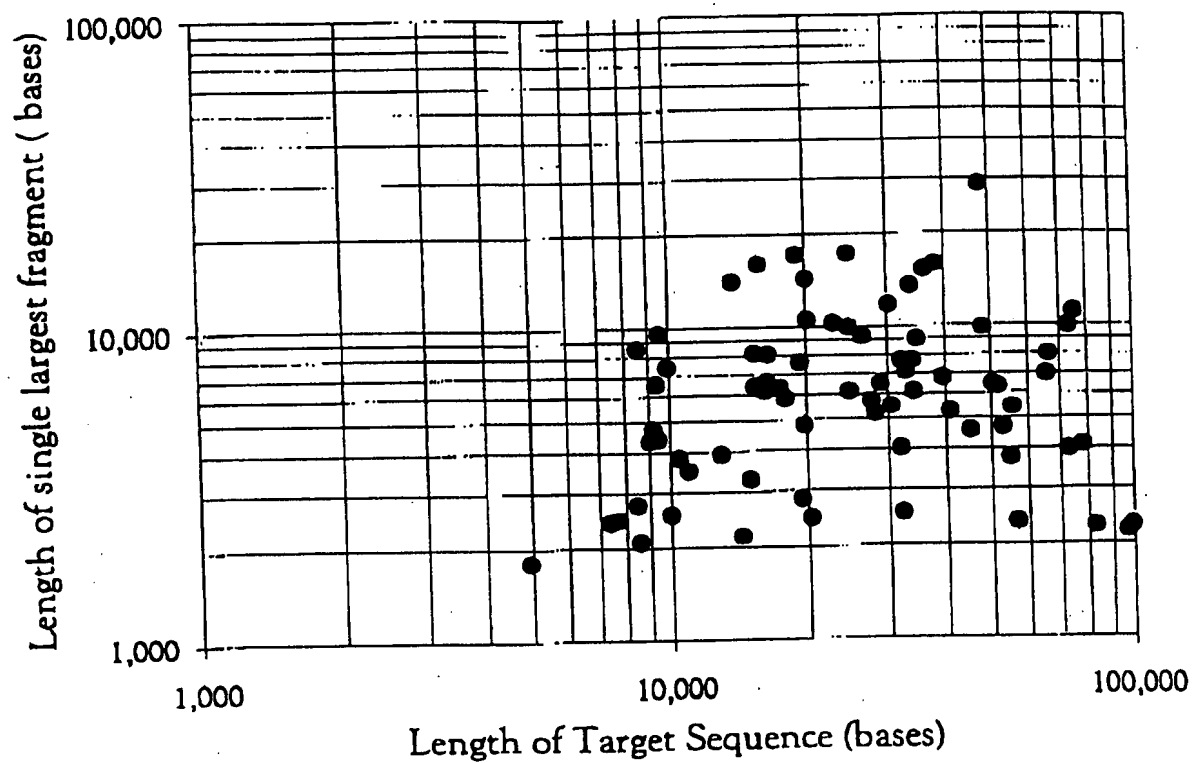


Figure 6

7/11

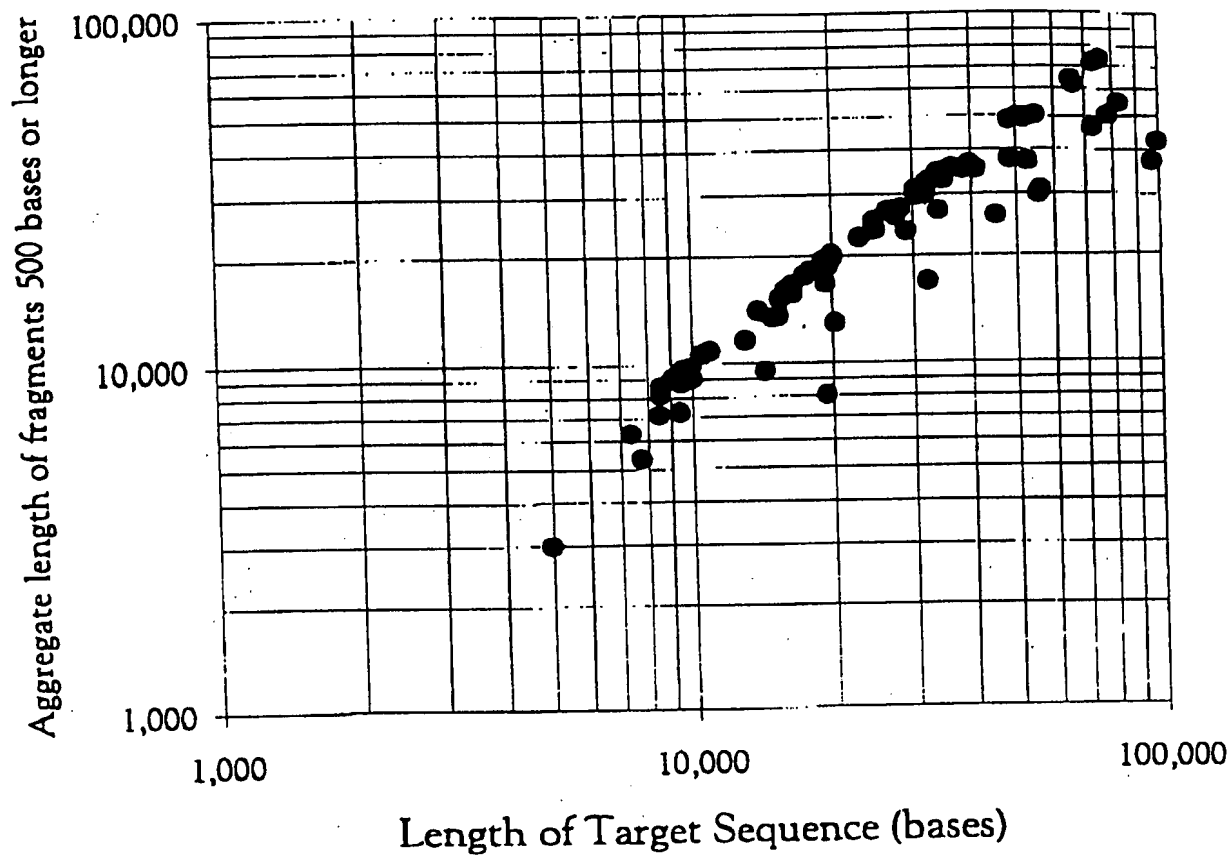


Figure 7

8/11

Local and Definition	Length	No. of different 16-mers in target	Fraction of 7-mers 16-mers that are repeated	No. of 16-mers in internal set	No. of 16-mers in internal set after error reduction	Fraction of internal set that are false positive	No. of fragments 500 bases or longer	Median length of fragments 500 bases or longer	Aggregate length of fragments 500 bases or longer	No. of fragments 100 bases or longer	Median length of fragments 100 bases or longer	Aggregate length of fragments 100 bases or longer	Trigram sequence coverage by fragments at least 100 bases in length (%)	Length of Ten Largest Fragments
[HARE-EURCP] Human oligodendrocyte myelin glycoprotein (OMG) exons 1-2; neurofascin-1 (NFI) exons 28-49; eukaryotic viral integration site 28 (EVIS28) exons 1-2; eukaryotic viral integration site 2A (EVISA) exons 1-2; eukaryotic virus (AVC) exons 1-2.	10049	9118	0.17	100735	18794	0.899	40	425	2584	272	272	7455	78.18	78.18
[HSE-KW364] Human PEA class III region containing cAMP response element binding protein-related protein (CREB-40) and intron 2 (intrac-2) genes, complete cds, complete sequence.	10023	5140	0.094	283252	48792	0.859	47	760	4208	155	155	6930	68.17	68.17
[F19623] Sequence of BAC F19623 from Arabidopsis thaliana chromosome 1, complete sequence.	69411	65804	0.027	1070700	96407	0.991	45	702	30743	243	207	79911	71.56	71.56
[HSA000009] Genomic sequence from Human 17, complete sequence.	8304	8174	0.035	343831	81177	0.810	57	844	54304	191	191	71111	68.13	68.13
[AC000400] Genomic sequence from Mouse 19, complete sequence.	76718	76417	0.002	1003018	75942	0.813	44	942	50045	130	130	6410	64.16	64.16
[U001477] 47 kb in 48 centomere region of Escherichia coli B402500.	75888	75810	0.001	1113669	75875	0.941	28	1880	77144	79	79	77157	95.31	95.31
[D10AC001646] Drosophila melanogaster (P1 D501845 (D071)) DNA sequence, complete sequence.	71727	71567	0.001	1003018	77942	0.931	71	801	71246	110	110	71315	91.33	91.33
[HARE-EUR] Human beta globin region on chromosome 11.	71308	71079	0.037	1003018	69713	0.877	41	460	46317	111	111	61302	61.30	61.30
[D10AC001653] Drosophila melanogaster (P1 D507876 (D14)) DNA sequence, complete sequence.	66930	66902	0.001	1003018	66708	0.877	79	1172	61559	56	56	66182	96.34	96.34
[SC051577] Saccharomyces cerevisiae chromosome V centromeres 9537, 9581, 9495, 9057, and centromere clone 509A.	66130	65870	0.002	75872	65716	0.913	27	27	1594	1776	1776	6594	97.36	97.36
[HARE-EUR] Human hypoxanthine phosphoribosyltransferase (HPR1) gene, complete cds.	80377	80311	0.003	166901	80360	0.877	33	744	10271	48	48	43580	78.77	78.77
[HSA0002117] Genomic sequence from Human 17, complete sequence.	55146	55117	0.001	1003018	55111	0.999	22	1003	55077	59	59	79152	79.15	79.15
[D10AC001655] Drosophila melanogaster (P1 D507700 (D171)) DNA sequence, complete sequence.	55146	55080	0.001	166901	54973	0.843	33	744	10271	48	48	43580	78.77	78.77
[BA001194] Enterobacter aerogenes plasmid R731, complete genome.	53103	53115	0.002	111596	52912	0.831	33	744	10271	48	48	43580	78.77	78.77
[HARE-EUR] Human hypoxanthine phosphoribosyltransferase (HPR1) gene, complete cds.	52737	52720	0.003	768644	52119	0.888	11	11	1153	1153	1153	52081	91.32	91.32
[SC051569] Saccharomyces cerevisiae chromosome V centromeres 9659, 9334, 9195, and centromere clone 1149.	50969	50811	0.003	118118	50638	0.645	24	24	1123	1123	1123	50969	91.32	91.32
[HARE-EUR] Human hypoxanthine phosphoribosyltransferase (HPR1) gene, complete cds.	48274	48231	0.002	166901	48193	0.871	17	17	17	17	17	48193	91.32	91.32
[LA000000] Saccharophaga thermophila, complete genome.	44502	44447	0.001	50312	44473	0.124	3	3	10460	10460	10460	44468	91.32	91.32

Figure 6 (1 of 4)

9/28/2007, EAST Version: 2.1.0.14

[illegible]

Figure 8 (2 of 4)

[illegible]

11/11

Local and Definition	Length	No. of different 16-mers in target	Fraction of Target 16-mers that are repeated	No. of 16-mers in internal set	No. of 16-mers in internal set after data reduction	Fraction of internal set that are the profiles	No. of fragments 500 bases or longer	Median Length of fragments 500 bases or longer	Approximate Length of fragments 500 bases or longer	No. of fragments 100 bases or longer	Median Length of fragments 100 bases or longer	Approximate Length of fragments 100 bases or longer	Target Sequence Coverage by fragments of length 100 bases or longer (%)	Length of Top Largest Fragment
[PUS94004] PUS measured (clone B100-17a, barcode-200) alone replication (P100P) gene, complete cds.	11031	10946	0.006	10943	10943	100.0	6	1414	10752	6	1414	10752	91.5	4059
[AF020102] Streptococcus pneumoniae G360 capsule polysaccharide locus, caps16f gene, partial cds, caps16d, caps16e, caps16f, caps16g, caps16h, caps16i and caps16k genes, complete cds, and caps16a, gene, partial cds.	10549	10549	0.000	10549	10549	100.0	3	3796	10449	3	3796	10449	98.03	2897
[U01000] Xenopus laevis genes for sperm-specific nuclear basic proteins (SP4), complete cds.	10143	7770	0.004	7287	6785	0.002	0	0	10752	0	0	10752	11.75	142
[D08195] Brucella abortus DNA for S-receptor kinase, complete cds.	10075	10075	0.001	10075	9953	0.002	4	1134	9822	9	755	9822	97.47	135
[AF020078] Arabidopsis thaliana kinase-like protein (P100-1) gene, complete cds.	10000	9944	0.000	9991	9876	0.001	2	4415	9822	3	222	9822	98.96	281
[AF020078] Arabidopsis thaliana acyl-CoA carboxylase gene, clone 1-11, complete cds.	9511	9511	0.000	9511	9442	0.000	1	1	9442	1	1	9442	98.96	947
[U01000] Clostridium magnum acyl gene, complete cds; TTP-dependent histone dehydrogenase alpha and beta subunit (acodA) genes, complete cds; acodB gene, dihydroisocitrate dehydrogenase (acodC) genes, complete cds.	9320	9320	0.016	9320	9320	0.000	4	1692	8732	4	1692	8732	91.9	104
[Z01000] Zee may's capsid-type retroelement PLEP-1 gag gene, complete cds.	9443	8151	0.135	8151	8107	0.000	1	1	8107	1	1	8107	82.32	182
[Z01000] Zee may's D2L, H (+)-transcribing ATPase (P101) gene, complete cds.	9235	9222	0.003	9224	9116	0.000	4	1894	9067	4	1894	9067	97.43	137
[H01000] Human immunodeficiency virus type 1 (HIV-1) proviral complete genome.	9113	9051	0.007	9053	8946	0.000	3	1264	8854	3	1264	8854	98.66	211
[U01000] Human histone H1B, H2A, H2B, and H4 genes, complete cds, and histone H3 gene, J and, gene cluster H3A1.	8608	8404	0.012	8403	8351	0.001	7	1030	8050	7	1030	8050	96.14	211
[H01000] Human mammary tumor virus proviral DNA (from P101) complete genome (P101) for glycoprotein and polyprotein and env protein, complete cds.	8501	8318	0.000	8318	8190	0.000	1	1	8190	1	1	8190	98.88	8505
[P01000] Fugu rubripes growth hormone (GH) gene, complete cds.	8211	7995	0.004	8000	7806	0.000	4	1400	7197	7	1179	7197	88.37	171
[U01000] Xenopus laevis 21-g trigger protein SPDL155 gene, complete cds.	7702	7772	0.000	7776	7659	0.000	4	1109	6510	11	555	6510	92.11	152
[U01000] Xenopus laevis elongation factor 1-alpha-O gene, complete cds.	7394	7348	0.004	7354	7234	0.001	4	1204	6510	7	7	6510	91.88	203
[U01000] Human lung carcinoma 10 kb secretory protein (CCO) gene, complete cds and Alu repeat sequences, complete cds.	4955	4884	0.000	4884	4725	0.001	3	589	4884	3	589	4884	90.33	110

Figure 8 (4 of 6)